

DECK 3A – BASIC MACHINE LEARNING APPLICATIONS

Dr. Jack Hong

Adjunct Faculty, Lee Kong Chian School of Business, SMU
Co-founder, Research Room Pte. Ltd.

INTRODUCTION TO MACHINE LEARNING

INTRODUCTION

- We will not cover all machine learning methods, but only those that are most effective and widely used
 - There are too many algorithms out there!
- Machine learning algorithms are classified into 2 categories:
 - Supervised learning
 - Unsupervised learning

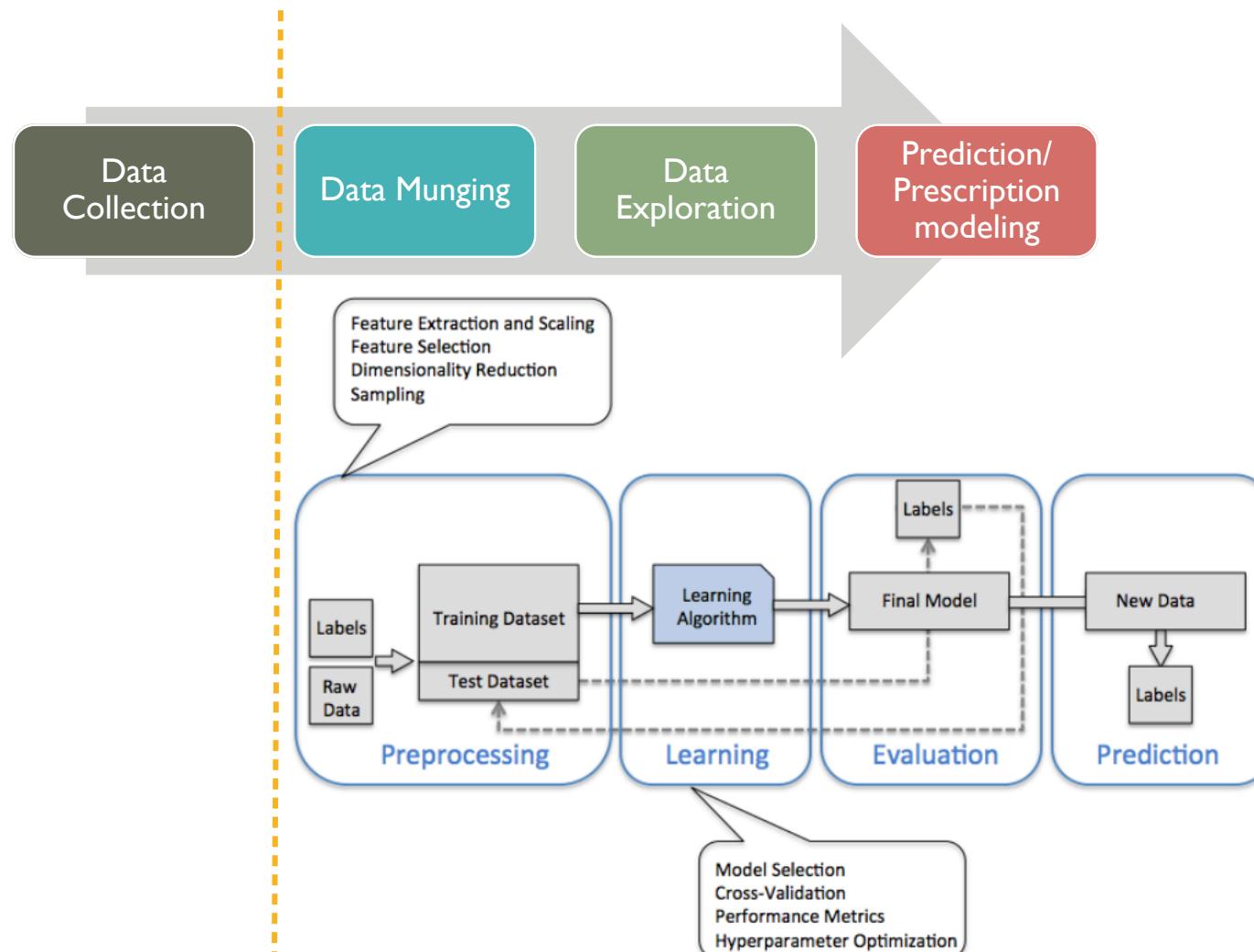
SUPERVISED LEARNING

- Train a model to learn from labeled training data
 - Each observation is classified with an outcome (e.g. Spam or Not Spam) or a continuous variable (as in regression)
 - We extract features from each observation
 - In the above example, the words that occur in each email
 - Can also include other variables such as time, name of sender, title etc
 - Compute a probability of outcome (spam or not spam) associated with the presence and absence of the features
 - In the case of the continuous variable, to minimize the loss function
- Reinforcement learning is about improving a model's performance based on its interaction with the environment
 - Uses reward function compared to loss function
 - E.g. Chess engine: System decides moves based on the state of the board, and the reward function will determine win or lose as the outcome

UNSUPERVISED LEARNING

- In supervised learning, we know the answer beforehand
- In unsupervised learning, we deal with unlabeled data
 - Unsupervised learning techniques allow us to explore the data structure to extract information without any prior guidance
 - Clustering: grouping observations by similarity and dissimilarity
 - Dimensionality reduction: removing noise from data while compressing high dimensional data into smaller ones

WHERE DOES MACHINE LEARNING FIT IN THE DATA SCIENCE VALUE CHAIN?



DATA MUNGING IN PYTHON

$$L(\mathbf{w}) = \prod_{i=1}^n P(y^i|x^i; \mathbf{w})$$

$$I_H(t) = -\sum_{i=1}^c p(i|t) \log_2 p(i|t)$$

$$k(x^i, x') = \exp\left(-\gamma \|x^i - x'\|^2\right)$$

$$p(Y \geq k) = \sum_{j=k}^n \binom{n}{j} \beta^j (1-\beta)^{n-j}$$

$$\varphi_2 = (u_2, v_2)$$

$$\varphi_3 = (u_3, v_3)$$

$$\varphi_4 = (u_4, v_4)$$

$$\varphi_5 = (u_5, v_5)$$

$$\varphi_6 = (u_6, v_6)$$

$$\varphi_7 = (u_7, v_7)$$

$$\varphi_8 = (u_8, v_8)$$

$$\varphi_9 = (u_9, v_9)$$

$$\varphi_{10} = (u_{10}, v_{10})$$

$$\varphi_{11} = (u_{11}, v_{11})$$

$$\varphi_{12} = (u_{12}, v_{12})$$

$$\varphi_{13} = (u_{13}, v_{13})$$

$$\varphi_{14} = (u_{14}, v_{14})$$

$$\varphi_{15} = (u_{15}, v_{15})$$

$$\varphi_{16} = (u_{16}, v_{16})$$

$$\varphi_{17} = (u_{17}, v_{17})$$

$$\varphi_{18} = (u_{18}, v_{18})$$

$$\varphi_{19} = (u_{19}, v_{19})$$

$$\varphi_{20} = (u_{20}, v_{20})$$

$$\varphi_{21} = (u_{21}, v_{21})$$

$$\varphi_{22} = (u_{22}, v_{22})$$

$$\varphi_{23} = (u_{23}, v_{23})$$

$$\varphi_{24} = (u_{24}, v_{24})$$

$$\varphi_{25} = (u_{25}, v_{25})$$

$$\varphi_{26} = (u_{26}, v_{26})$$

$$\varphi_{27} = (u_{27}, v_{27})$$

$$\varphi_{28} = (u_{28}, v_{28})$$

$$\varphi_{29} = (u_{29}, v_{29})$$

$$\varphi_{30} = (u_{30}, v_{30})$$

$$\varphi_{31} = (u_{31}, v_{31})$$

$$\varphi_{32} = (u_{32}, v_{32})$$

$$\varphi_{33} = (u_{33}, v_{33})$$

$$\varphi_{34} = (u_{34}, v_{34})$$

$$\varphi_{35} = (u_{35}, v_{35})$$

$$\varphi_{36} = (u_{36}, v_{36})$$

$$\varphi_{37} = (u_{37}, v_{37})$$

$$\varphi_{38} = (u_{38}, v_{38})$$

$$\varphi_{39} = (u_{39}, v_{39})$$

$$\varphi_{40} = (u_{40}, v_{40})$$

$$\varphi_{41} = (u_{41}, v_{41})$$

$$\varphi_{42} = (u_{42}, v_{42})$$

$$\varphi_{43} = (u_{43}, v_{43})$$

$$\varphi_{44} = (u_{44}, v_{44})$$

$$\varphi_{45} = (u_{45}, v_{45})$$

$$\varphi_{46} = (u_{46}, v_{46})$$

$$\varphi_{47} = (u_{47}, v_{47})$$

$$\varphi_{48} = (u_{48}, v_{48})$$

$$\varphi_{49} = (u_{49}, v_{49})$$

$$\varphi_{50} = (u_{50}, v_{50})$$

$$\varphi_{51} = (u_{51}, v_{51})$$

$$\varphi_{52} = (u_{52}, v_{52})$$

$$\varphi_{53} = (u_{53}, v_{53})$$

$$\varphi_{54} = (u_{54}, v_{54})$$

$$\varphi_{55} = (u_{55}, v_{55})$$

$$\varphi_{56} = (u_{56}, v_{56})$$

$$\varphi_{57} = (u_{57}, v_{57})$$

$$\varphi_{58} = (u_{58}, v_{58})$$

$$\varphi_{59} = (u_{59}, v_{59})$$

$$\varphi_{60} = (u_{60}, v_{60})$$

$$\varphi_{61} = (u_{61}, v_{61})$$

$$\varphi_{62} = (u_{62}, v_{62})$$

$$\varphi_{63} = (u_{63}, v_{63})$$

$$\varphi_{64} = (u_{64}, v_{64})$$

$$\varphi_{65} = (u_{65}, v_{65})$$

$$\varphi_{66} = (u_{66}, v_{66})$$

$$\varphi_{67} = (u_{67}, v_{67})$$

$$\varphi_{68} = (u_{68}, v_{68})$$

$$\varphi_{69} = (u_{69}, v_{69})$$

$$\varphi_{70} = (u_{70}, v_{70})$$

$$\varphi_{71} = (u_{71}, v_{71})$$

$$\varphi_{72} = (u_{72}, v_{72})$$

$$\varphi_{73} = (u_{73}, v_{73})$$

$$\varphi_{74} = (u_{74}, v_{74})$$

$$\varphi_{75} = (u_{75}, v_{75})$$

$$\varphi_{76} = (u_{76}, v_{76})$$

$$\varphi_{77} = (u_{77}, v_{77})$$

$$\varphi_{78} = (u_{78}, v_{78})$$

$$\varphi_{79} = (u_{79}, v_{79})$$

$$\varphi_{80} = (u_{80}, v_{80})$$

$$\varphi_{81} = (u_{81}, v_{81})$$

$$\varphi_{82} = (u_{82}, v_{82})$$

$$\varphi_{83} = (u_{83}, v_{83})$$

$$\varphi_{84} = (u_{84}, v_{84})$$

$$\varphi_{85} = (u_{85}, v_{85})$$

$$\varphi_{86} = (u_{86}, v_{86})$$

$$\varphi_{87} = (u_{87}, v_{87})$$

$$\varphi_{88} = (u_{88}, v_{88})$$

$$\varphi_{89} = (u_{89}, v_{89})$$

$$\varphi_{90} = (u_{90}, v_{90})$$

$$\varphi_{91} = (u_{91}, v_{91})$$

$$\varphi_{92} = (u_{92}, v_{92})$$

$$\varphi_{93} = (u_{93}, v_{93})$$

$$\varphi_{94} = (u_{94}, v_{94})$$

$$\varphi_{95} = (u_{95}, v_{95})$$

$$\varphi_{96} = (u_{96}, v_{96})$$

$$\varphi_{97} = (u_{97}, v_{97})$$

$$\varphi_{98} = (u_{98}, v_{98})$$

$$\varphi_{99} = (u_{99}, v_{99})$$

$$\varphi_{100} = (u_{100}, v_{100})$$

$$\varphi_{101} = (u_{101}, v_{101})$$

$$\varphi_{102} = (u_{102}, v_{102})$$

$$\varphi_{103} = (u_{103}, v_{103})$$

$$\varphi_{104} = (u_{104}, v_{104})$$

$$\varphi_{105} = (u_{105}, v_{105})$$

$$\varphi_{106} = (u_{106}, v_{106})$$

$$\varphi_{107} = (u_{107}, v_{107})$$

$$\varphi_{108} = (u_{108}, v_{108})$$

$$\varphi_{109} = (u_{109}, v_{109})$$

$$\varphi_{110} = (u_{110}, v_{110})$$

$$\varphi_{111} = (u_{111}, v_{111})$$

$$\varphi_{112} = (u_{112}, v_{112})$$

$$\varphi_{113} = (u_{113}, v_{113})$$

$$\varphi_{114} = (u_{114}, v_{114})$$

$$\varphi_{115} = (u_{115}, v_{115})$$

$$\varphi_{116} = (u_{116}, v_{116})$$

$$\varphi_{117} = (u_{117}, v_{117})$$

$$\varphi_{118} = (u_{118}, v_{118})$$

$$\varphi_{119} = (u_{119}, v_{119})$$

$$\varphi_{120} = (u_{120}, v_{120})$$

$$\varphi_{121} = (u_{121}, v_{121})$$

$$\varphi_{122} = (u_{122}, v_{122})$$

$$\varphi_{123} = (u_{123}, v_{123})$$

$$\varphi_{124} = (u_{124}, v_{124})$$

$$\varphi_{125} = (u_{125}, v_{125})$$

$$\varphi_{126} = (u_{126}, v_{126})$$

$$\varphi_{127} = (u_{127}, v_{127})$$

$$\varphi_{128} = (u_{128}, v_{128})$$

$$\varphi_{129} = (u_{129}, v_{129})$$

$$\varphi_{130} = (u_{130}, v_{130})$$

$$\varphi_{131} = (u_{131}, v_{131})$$

$$\varphi_{132} = (u_{132}, v_{132})$$

$$\varphi_{133} = (u_{133}, v_{133})$$

$$\varphi_{134} = (u_{134}, v_{134})$$

$$\varphi_{135} = (u_{135}, v_{135})$$

$$\varphi_{136} = (u_{136}, v_{136})$$

$$\varphi_{137} = (u_{137}, v_{137})$$

$$\varphi_{138} = (u_{138}, v_{138})$$

$$\varphi_{139} = (u_{139}, v_{139})$$

$$\varphi_{140} = (u_{140}, v_{140})$$

$$\varphi_{141} = (u_{141}, v_{141})$$

$$\varphi_{142} = (u_{142}, v_{142})$$

$$\varphi_{143} = (u_{143}, v_{143})$$

$$\varphi_{144} = (u_{144}, v_{144})$$

$$\varphi_{145} = (u_{145}, v_{145})$$

$$\varphi_{146} = (u_{146}, v_{146})$$

$$\varphi_{147} = (u_{147}, v_{147})$$

$$\varphi_{148} = (u_{148}, v_{148})$$

$$\varphi_{149} = (u_{149}, v_{149})$$

$$\varphi_{150} = (u_{150}, v_{150})$$

$$\varphi_{151} = (u_{151}, v_{151})$$

$$\varphi_{152} = (u_{152}, v_{152})$$

$$\varphi_{153} = (u_{153}, v_{153})$$

$$\varphi_{154} = (u_{154}, v_{154})$$

$$\varphi_{155} = (u_{155}, v_{155})$$

$$\varphi_{156} = (u_{156}, v_{156})$$

$$\varphi_{157} = (u_{157}, v_{157})$$

$$\varphi_{158} = (u_{158}, v_{158})$$

$$\varphi_{159} = (u_{159}, v_{159})$$

$$\varphi_{160} = (u_{160}, v_{160})$$

$$\varphi_{161} = (u_{161}, v_{161})$$

$$\varphi_{162} = (u_{162}, v_{162})$$

$$\varphi_{163} = (u_{163}, v_{163})$$

$$\varphi_{164} = (u_{164}, v_{164})$$

$$\varphi_{165} = (u_{165}, v_{165})$$

$$\varphi_{166} = (u_{166}, v_{166})$$

$$\varphi_{167} = (u_{167}, v_{167})$$

$$\varphi_{168} = (u_{168}, v_{168})$$

$$\varphi_{169} = (u_{169}, v_{169})$$

$$\varphi_{170} = (u_{170}, v_{170})$$

$$\varphi_{171} = (u_{171}, v_{171})$$

$$\varphi_{172} = (u_{172}, v_{172})$$

$$\varphi_{173} = (u_{173}, v_{173})$$

$$\varphi_{174} = (u_{174}, v_{174})$$

$$\varphi_{175} = (u_{175}, v_{175})$$

$$\varphi_{176} = (u_{176}, v_{176})$$

$$\varphi_{177} = (u_{177}, v_{177})$$

$$\varphi_{178} = (u_{178}, v_{178})$$

$$\varphi_{179} = (u_{179}, v_{179})$$

$$\varphi_{180} = (u_{180}, v_{180})$$

$$\varphi_{181} = (u_{181}, v_{181})$$

$$\varphi_{182} = (u_{182}, v_{182})$$

$$\varphi_{183} = (u_{183}, v_{183})$$

$$\varphi_{184} = (u_{184}, v_{184})$$

$$\varphi_{185} = (u_{185}, v_{185})$$

$$\varphi_{186} = (u_{186}, v_{186})$$

$$\varphi_{187} = (u_{187}, v_{187})$$

$$\varphi_{188} = (u_{188}, v_{188})$$

$$\varphi_{189} = (u_{189}, v_{189})$$

$$\varphi_{190} = (u_{190}, v_{190})$$

$$\varphi_{191} = (u_{191}, v_{191})$$

$$\varphi_{192} = (u_{192}, v_{192})$$

$$\varphi_{193} = (u_{193}, v_{193})$$

$$\varphi_{194} = (u_{194}, v_{194})$$

$$\varphi_{195} = (u_{195}, v_{195})$$

$$\varphi_{196} = (u_{196}, v_{196})$$

$$\varphi_{197} = (u_{197}, v_{197})$$

$$\varphi_{198} = (u_{198}, v_{198})$$

$$\varphi_{199} = (u_{$$

DATA MUNGING IN PYTHON

- Data munging is a necessary but dreaded process in data analytics
 - Involves cleaning, transforming, and setting data the right way
- Compared to R, Python is more powerful (and convenient) in data munging
 - Data scientists use the pandas package for this purpose
 - Pandas is developed by Wes McKinney, a Financial Economist
 - a financial economist to handle economics and finance datasets (refer to panel data setup in Lesson 3)
 - Pandas got its name not from the animal but because it was meant for handling **panel data**.
 - The package was so useful that its now the de-facto data munging package in Python

IMPORTANT TECHNIQUES USING PANDAS AND SCIKIT-LEARN

- I'll cover a few key techniques that you will use frequently
 - Upload “ML introduction.ipynb” and “wine.csv” into your Jupyter dashboard
 - Follow the comments and instructions
 - This notebook tutorial covers data munging and predictive modeling
- To know more about pandas, you may refer to <https://pandas.pydata.org/pandas-docs/stable/10min.html> for the official 10 minutes guide to pandas

PERCEPTRON LEARNING

WHERE IT ALL STARTED FOR MACHINE
LEARNING

$$L(\mathbf{w}) = \prod_{i=1}^n P(y^i|x^i; \mathbf{w})$$

$$I_H(t) = -\sum_{i=1}^2 p(i|t) \log_2 p(i|t)$$

$$k(x^i, x') = \exp(-\gamma \|x^i - x'\|^2)$$

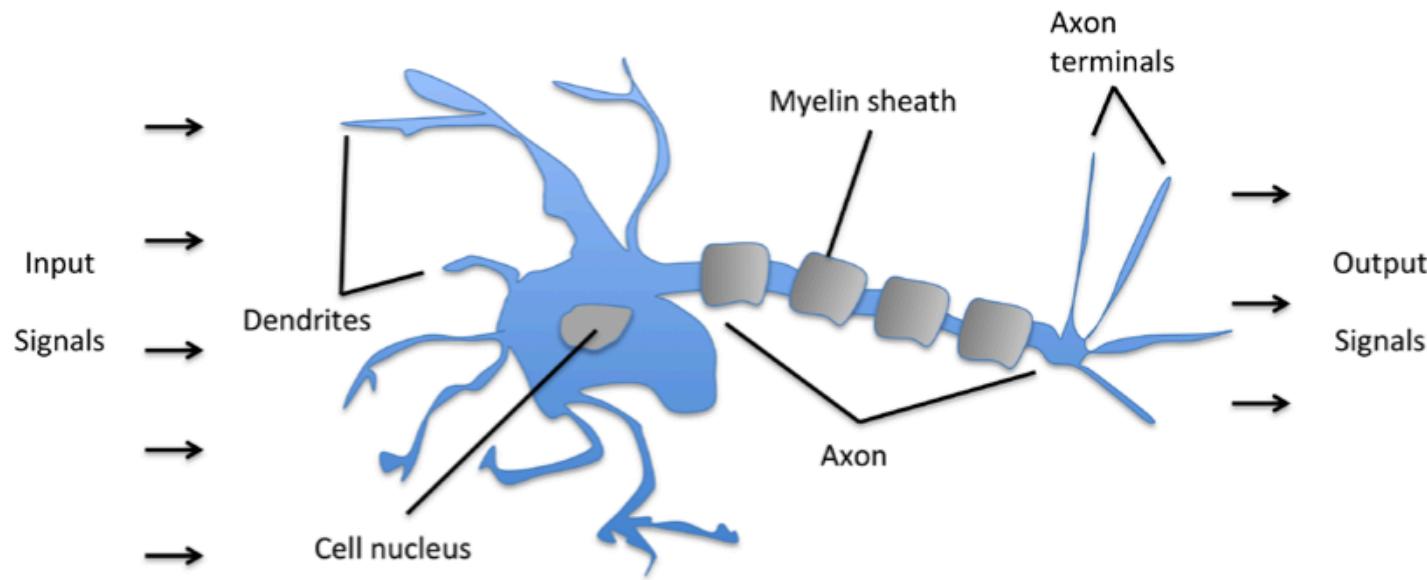
$$p(Y \geq k) = \sum_{j=k}^n \binom{n}{j} \delta_j^{k-1}$$

$$\varphi_y = \eta_y \varphi_z(z, \varphi_x)$$

$$z_i = \phi(\varphi_x \cdot \varphi_y \cdot \varphi_z)$$

THE BEGINNINGS OF MACHINE LEARNING

- Early works in machine learning attempts to mimic how a brain cell works, in simplistic forms
 - Simple logic gate with binary outputs, triggered by the level/strength of accumulated signal against a threshold

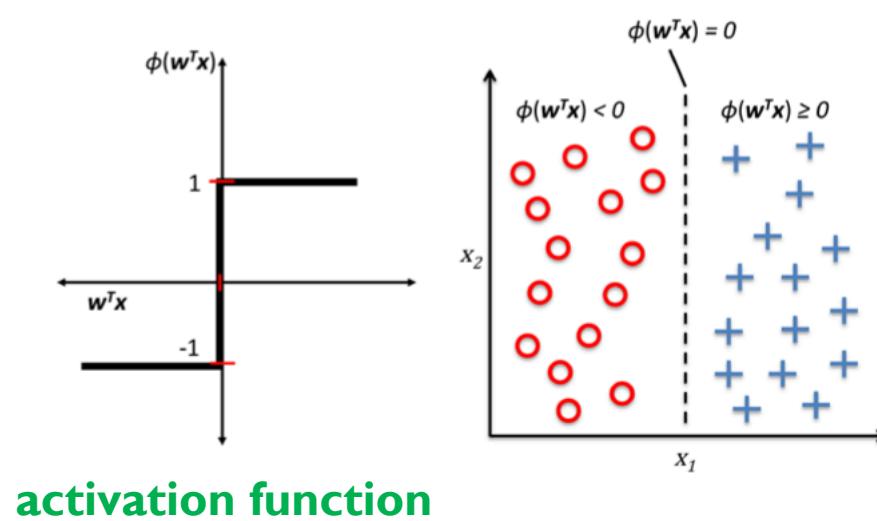


PERCEPTRON LEARNING – BASIC ML ILLUSTRATION

- Based on the Neuron model, Frank Rosenblatt conceptualized a learning rule that would
 - Automatically learn the optimal weight coefficients
 - Multiplied with the input features
 - Decide whether a neuron fires or not
- This is akin to a classification task
- In the below example, we use a step function (if $x > 0.5$, $y = \text{red}$) (else $y = \text{blue}$)
 - All values of x will then be mapped into a y -value of 1 or 0
 - Example: If humidity is $> 60\%$, it will rain, else it will not

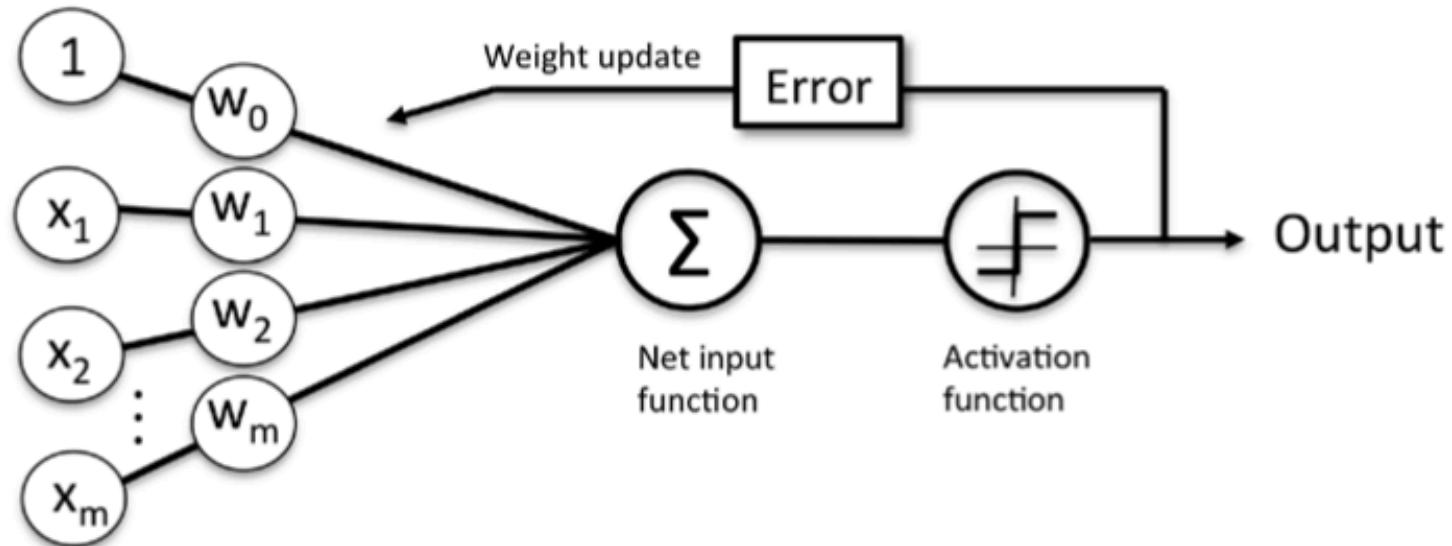
$$\mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

weights **input**



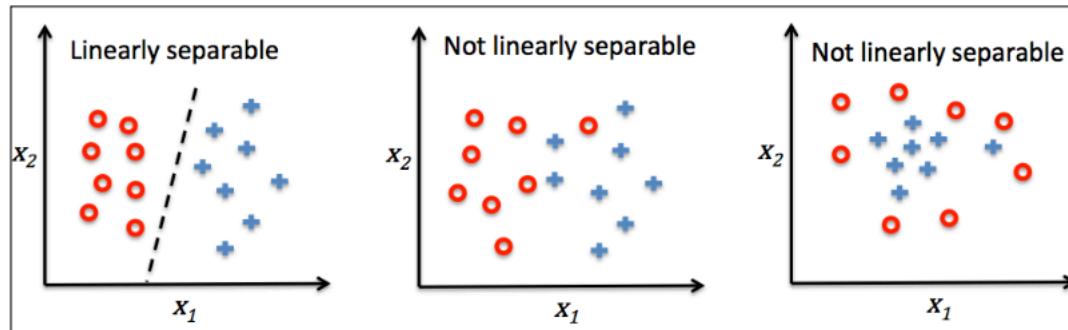
HOW DOES THE MODEL LEARN?

- An ML model learns by iteration
- It chooses a set of starting weights for each feature input
- Throws it through the activation function and compares the predicted output against the actual
- The model continuously attempts to reduce the error (actual – predicted) by changing the weights for the input features



PERCEPTRON LEARNING

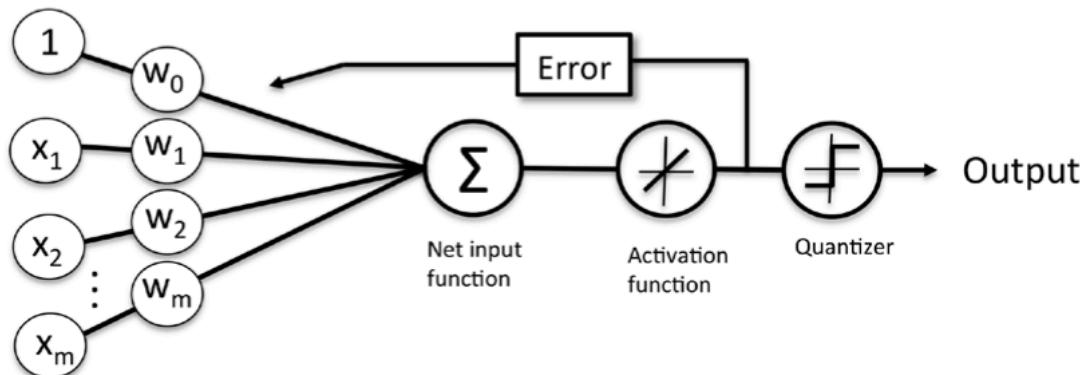
- Intuitively, we can see that this method adopts a linear separation mechanism to split samples via a boundary
 - However, if the samples cannot be perfectly separated, then the perceptron will never terminate
 - The solution is to allow some kind of error threshold (e.g. 20% misclassification)



- We can generalize this perceptron model and change the step function to other functional forms
 - Such as the linear and logistic models we learnt in Lesson 3
- Machine learning is very closely related to statistical modeling

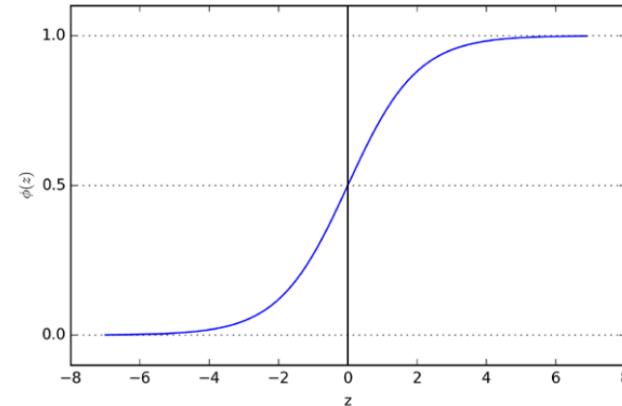
LINEAR REGRESSION IN MACHINE LEARNING

- The linear regression model is implemented in machine learning by using a linear activation function
- This ML model is known as Adaline

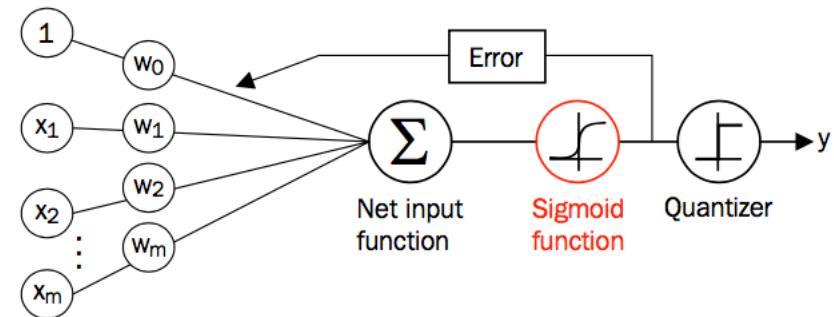


LOGISTIC REGRESSION IN MACHINE LEARNING

- Recall that logistic regression uses a logit function to map values of x to a range bounded by $(0, 1)$
- In machine learning, such a function is known as the sigmoid function due to its S-shape
- Implementing the logistic regression in ML is simply using the sigmoid function as the activation function



Sigmoid function



GRADIENT DESCENT I

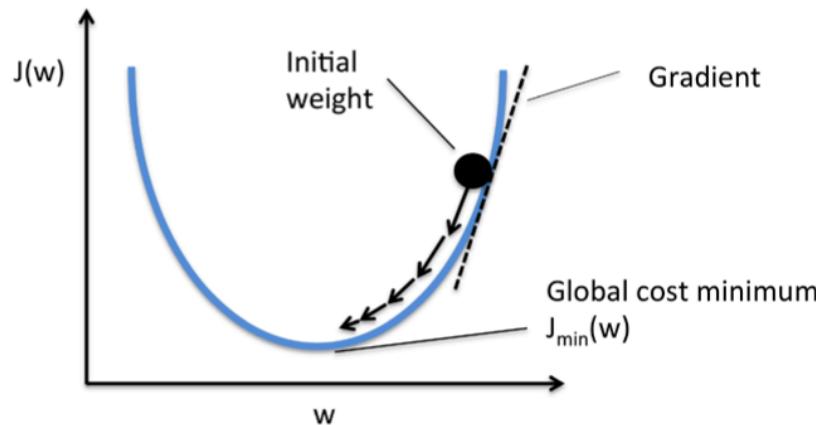
- By changing the activation function to a linear function, we are attempting to fit data into a linear form, similar to the linear regression
 - The model error is now computed based on a linear prediction
- Like regression models, the objective of the ML algorithm is to find a set of weights for each input feature that minimizes predictive errors (mathematically, it's defined as the sum of squared errors)

$$J(\mathbf{w}) = \frac{1}{2} \sum_i (y^{(i)} - \phi(z^{(i)}))^2$$

- This is the general principle of supervised machine learning algorithms:
 - Define an objective function to be optimized by reducing the error between predicted and actual values
 - Shift the weights to find the minimum error rate achievable (in machine learning speak, we also call error the cost)

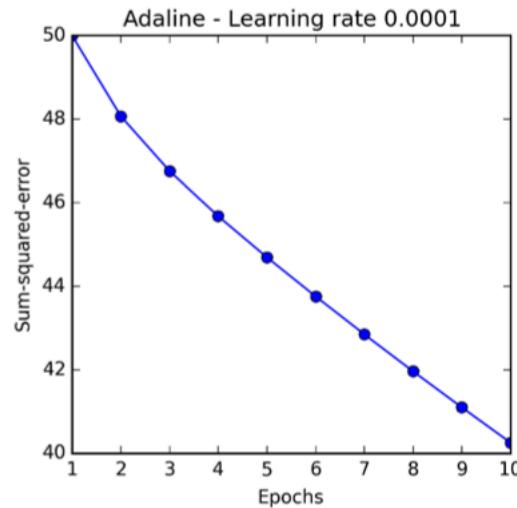
GRADIENT DESCENT II

- Analogous to climbing down a hill to reach the bottom of the valley

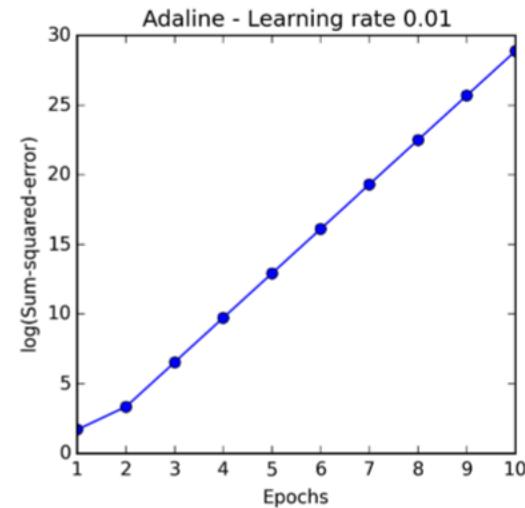


- Hyper-parameters
 - For the above model, there are 2 parameters that we need to care about as we perform the gradient descent
 - Number of iterations: How many times should we move the weights around?
 - Learning rate: How much movement should each iteration allow?

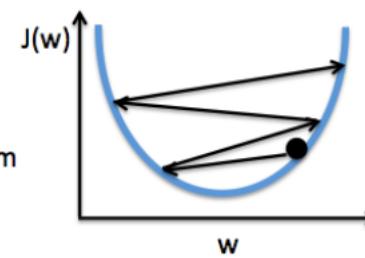
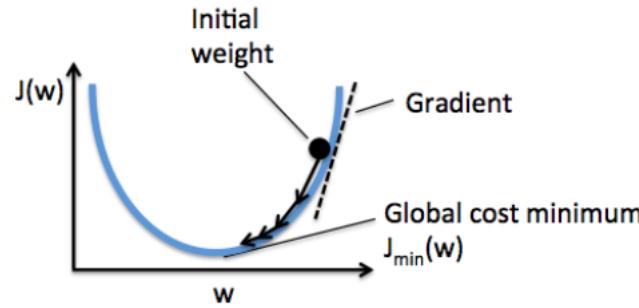
CHOOSING HYPERPARAMETERS



A small learning rate requires a large amount of iterations to converge



However, if the learning rate is too large, we keep overshooting the global minimum and errors increase



- Hyperparameter tuning (finding the best hyperparameter) is an indispensable step in machine learning modeling
 - There are many tricks to hyperparameter tuning but we will not cover them in this course

PREDICTIVE ACCURACY

$$L(\mathbf{w}) = \prod_{i=1}^n P(y^i|x^i; \mathbf{w}) \quad I_H(t) = -\sum_{i=1}^c p(i|t) \log_2 p(i|t) \quad k(x^i, x') = \exp(-\gamma \|x^i - x'\|^2) \quad p(Y \geq k) = \sum_{j=k}^n \binom{n}{j} \beta^j (1-\beta)^{n-j}$$

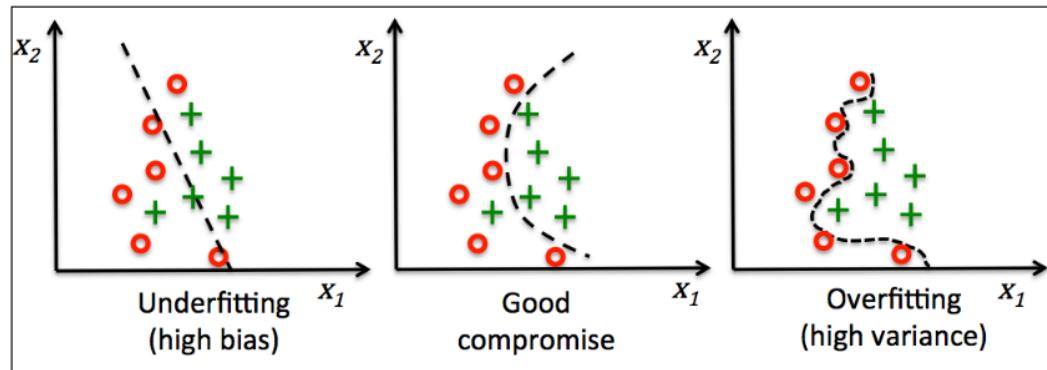
$$\mathcal{Z}[\mathcal{I}_1], \mathcal{Z}[\mathcal{J}_1]$$

$$size_{\leq i} \phi e^{\lambda_3 X} \lambda_3^3$$

$$\varphi_2 = (u_2, v_2)$$

OVERFITTING AND UNDERFITTING

- Recall in Lesson 3, we had to split the dataset into training and test set
 - Because a model tends to perform well on data it has “seen before” (in-sample data)
 - This has an upward bias on reported predictive accuracy
 - Performs well on training data, but badly on out-sample data
 - The solution is to train the model in-sample and test it on out-sample data
 - The name of this problem is known as “overfitting”
- In contrast to overfitting, underfitting is another issue where the model is too naïve (e.g. too little features) to capture significant patterns in the training data
 - Reported predictive accuracy is weak for both in and out-sample data



REGULARIZATION TO COMBAT OVER AND UNDER FITTING

- Finding a good trade-off requires some form of tuning on the complexity of the model
- Regularization is a good method for this purpose
 - Can handle collinearity and filter noise
 - Introduce bias to penalize extreme parameter weights
 - Called L1 or L2 regularization, or a mix of both
 - Most commonly used: L2 regularization (also called shrinkage or weight decay)
- Regularization parameters are included as hyperparameters in most ML algorithms in Python's scikit-learn package

SUPPORT VECTOR MACHINES

$$L(\mathbf{w}) = \prod_{i=1}^n P(y^i|x^i; \mathbf{w}) \quad I_H(t) = -\sum_{i=1}^c p(i|t) \log_2 p(i|t) \quad k(x^i, x^j) = \exp(-\gamma \|x^i - x^j\|^2) \quad p(Y \geq k) = \sum_{y=k}^n \binom{n}{y} \beta^y (1-\beta)^{n-y}$$

92

=

(n_r, z_r)

x

y_r

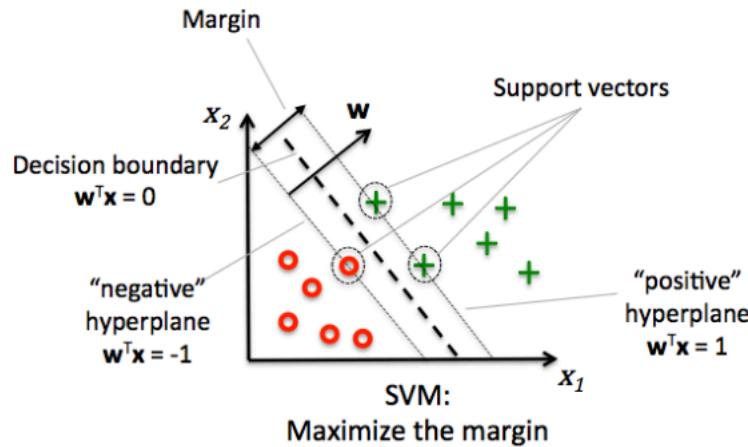
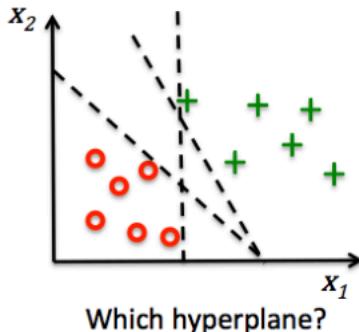
φ_r

z_r

z<sub

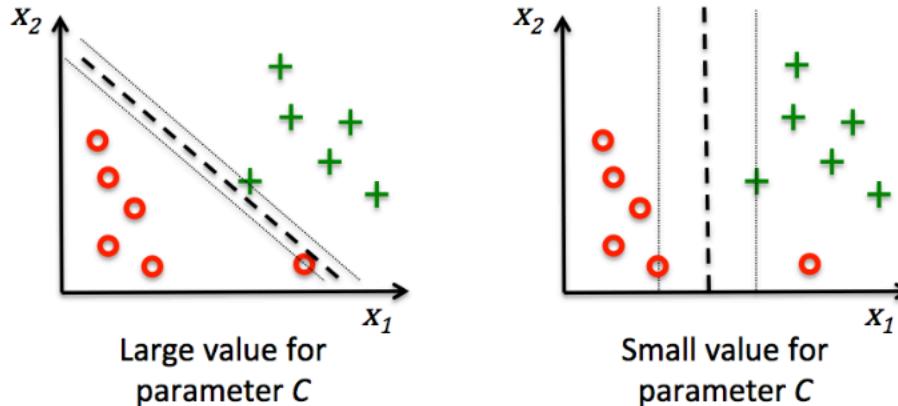
SUPPORT VECTOR MACHINES (SVM)

- A popular and powerful ML algorithm extended from the Perceptron model
- Recall the optimization objective of the Perceptron model is to minimize errors
 - SVM's optimization objective is to maximize the distance between observations separated by a decision boundary called the hyperplane
 - In other words, SVM optimizes by maximizing the margin
 - Note that hyperplanes can be multi-dimensional (which is often the case)



SEPARATION WITH ERRORS

- In the event that a perfect cut cannot be established, the SVM model allows for a threshold of misclassification errors
 - In the below example, the parameter C defines this threshold



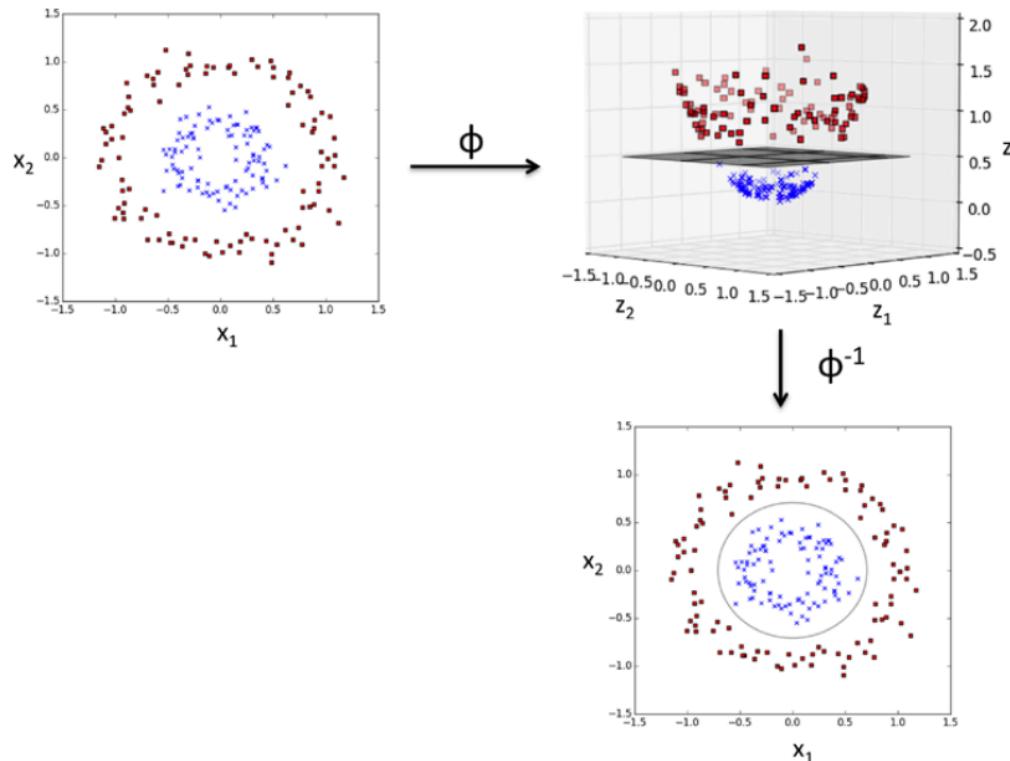
The value of C corresponds to the size of the penalties for misclassification errors

SVM FOR NONLINEAR PROBLEMS

- Kernel SVM
 - Used when its not possible for us to separate positive and negative outcomes using linear methods
 - The idea is to create nonlinear combinations of the original features/variables and project them onto a higher dimensional space via a mapping function
 - then we can do a linear hyperplane in this higher dimensional space
 - map it back to the original lower dimensional space after we are done

MAPPING FUNCTION

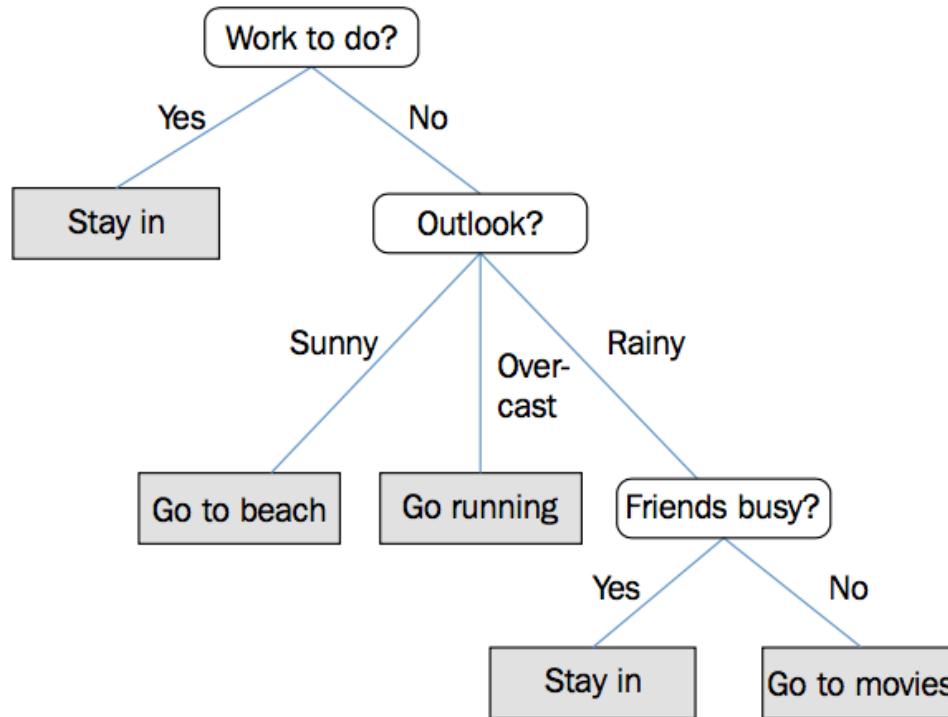
- Example:
 - Transform a 2-dimensional dataset into a 3-dimensional one
 - In this case, the closer the values of X_1 and X_2 are to the center, the lower the Z value is
 - Use a linear function to cut a line between the red and blue observations
 - Transform the results back into 2-d



DECISION TREES

DECISION TREES

- Decision trees are so named because the model breaks the data down by decision-based questions
 - Features can be categorical, as well as continuous



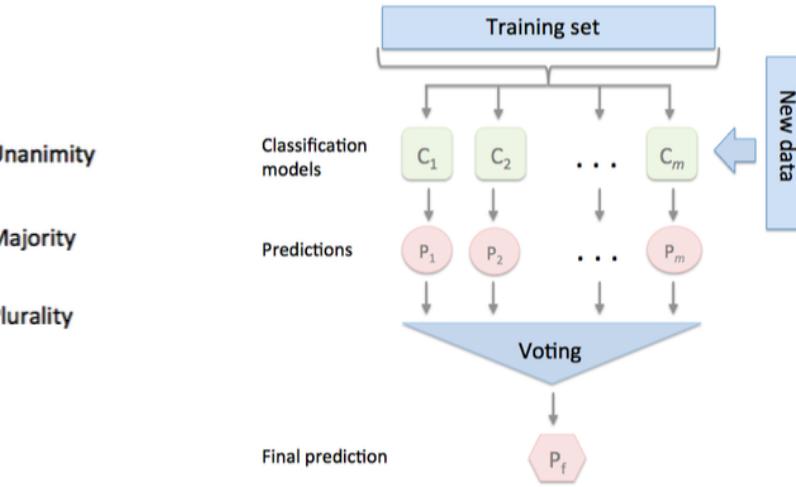
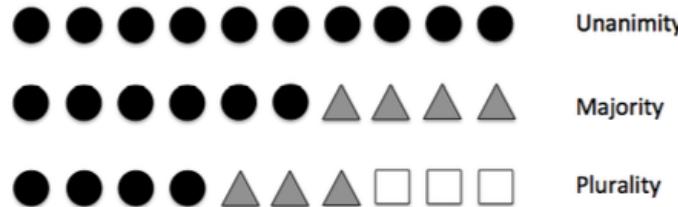
IMPLEMENTING DECISION TREE LEARNING

- The splitting rule at each node of the tree acts on the feature that results in the largest Information Gain (IG)
 - This is repeated at each child node until the observations at each node all belong to the same class
- As such, decision tree models are extremely good at fitting
 - The model can create, shrink, and expand boxes to envelop datapoints with the same outcome variable
 - However, such a procedure is prone to overfitting
 - Techniques such as pruning (sets a limit for the maximal depth of the tree) help to mitigate this

ENSEMBLE METHODS

WHAT ARE ENSEMBLE METHODS?

- Combining individual classifiers into a meta-classifier can realize better out-of-sample performance
- Generally, ensemble methods collect predictions from multiple models and use a voting mechanism to choose the best result



RANDOM FOREST

- Random Forest is a type of ensemble method built on decision trees
 - RF is simply a collection of decision trees
- A very simple adaptation but often results in drastic improvements to predictive accuracy
 - Good performance, scalable and easy to use
- Combining weak learners to build a strong learner
 - The idea that multiple independent models give better results than a single model
 - Assuming 100 classifiers that are each correct only 55% of the time
 - However, the probability that the majority of them are correct increases to 82% (cumulative binomial probability)
- Other ensemble methods: bagging and boosting

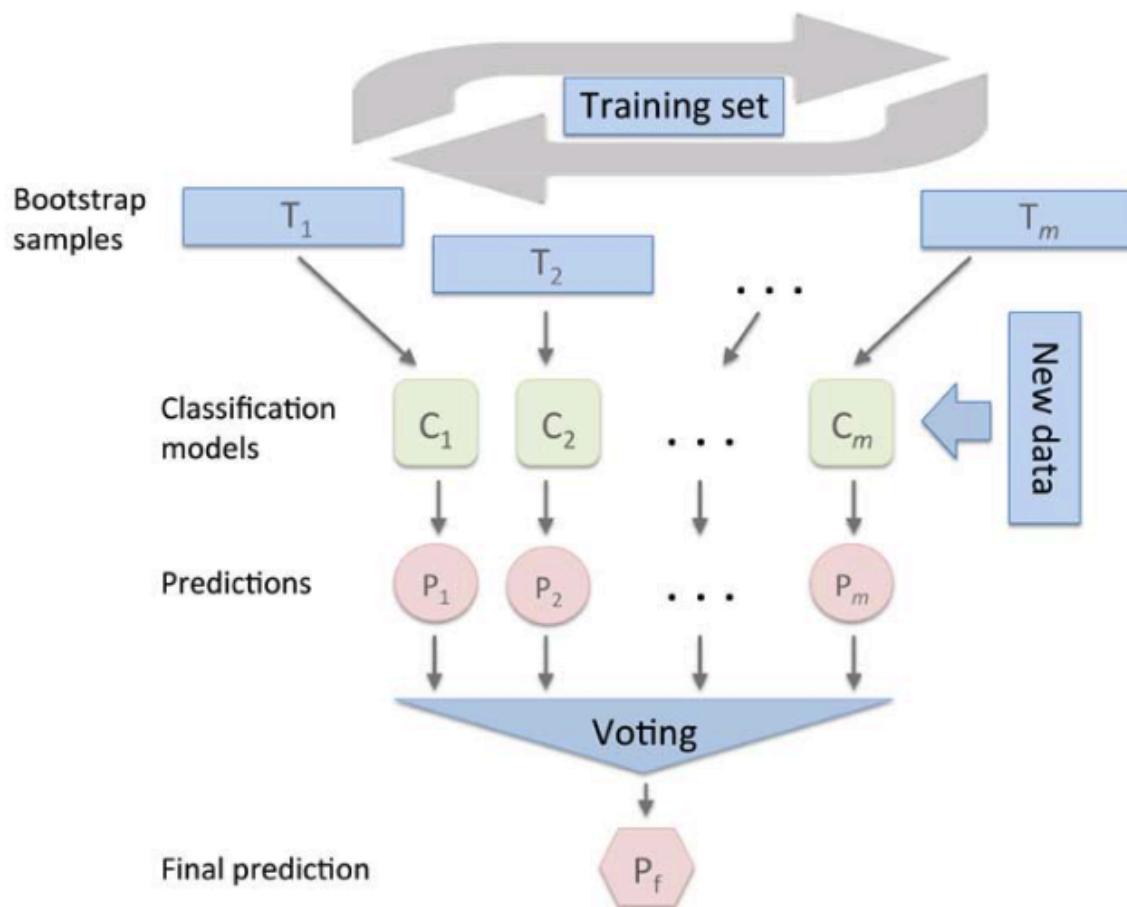
IMPLEMENTING RANDOM FOREST

- Bootstrap a sample of size n
 - Bootstrapping is a method of sample selection where we make random selections with replacement
- Grow a decision tree using this sample
 - At each node, randomly select d features without replacement
 - Split this node using the best feature
- Repeat the above steps k times
- Aggregate the prediction by each tree and assign the class to the data via majority voting

SOME REMARKS

- Random forest model is a black box
 - Cannot be interpreted like a decision tree where you know that branches come out of important features
 - Random forest is a bunch of trees, and the prediction is a result of majority voting (i.e. most of the trees predict this outcome)
- However, RF has the following advantages
 - No need to worry about hyperparameter values
 - Don't need to prune the forest
 - Only need to choose the number of trees
 - The larger the number, the better the performance
 - However, the larger the number, the more expensive the computational cost

BAGGING – AN ALTERNATIVE TO MAJORITY VOTING

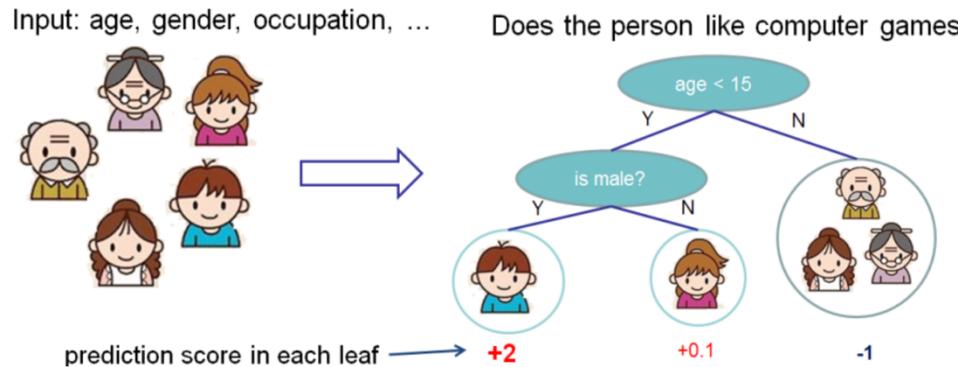


BOOSTING

- 3 commonly-used techniques
 - Adaptive boosting (adaboost)
 - A type of classifier boosting
 - Stochastic Gradient boosting
 - A type of regression boosting
 - Extreme gradient boosting (xgboost)
 - A type of tree ensemble (combination of classification and regression trees)
 - Fits a regularization variable in the objective function to prevent over-fitting
- The basic concept behind boosting is to start with very simple base classifiers (slightly better than random guessing)
 - These simple classifiers are called weak learners (e.g. 2-layer decision tree – they call this a stump)
 - Continuously make the weak learners learn from misclassified samples
 - Draw a random subset without replacement and train a weak learner
 - Draw a second subset without replacement, add 50% of the misclassified samples in (1), and train another weak learner
 - Find the training samples that both (1) and (2) disagree, and train a third weak learner
 - Combine all 3 learners via majority voting
 - Note that boosting algorithms have a tendency to overfit

XGBOOST

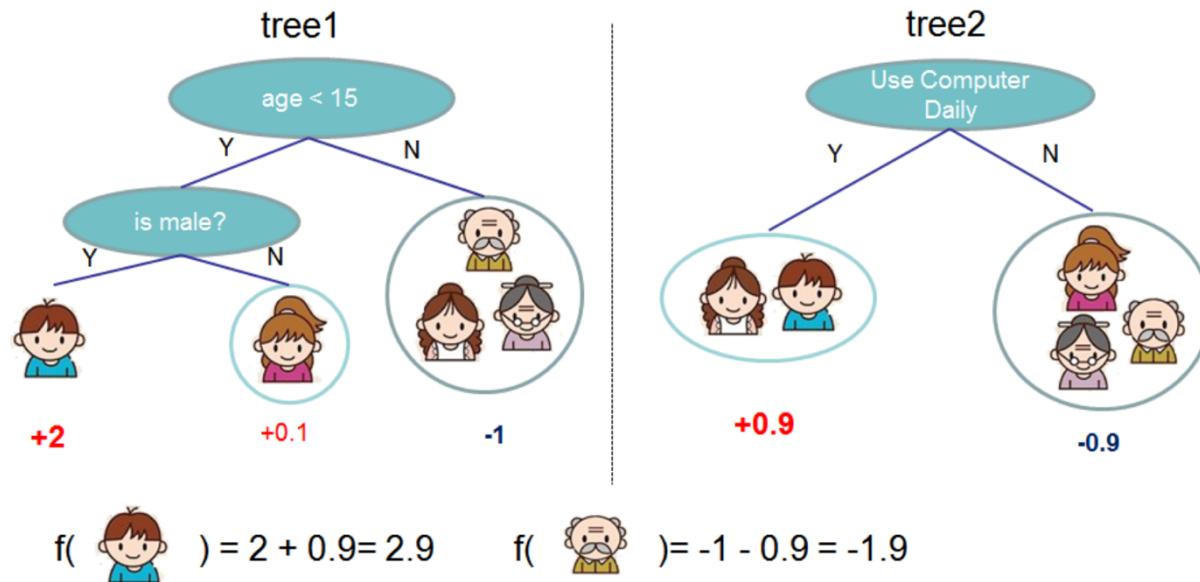
- Tree ensemble
 - Combination of classification and regression trees (CART)



- In normal decision trees, the leaf contains the decision values
 - In CART, the leaves contain a real-valued score

TREE ENSEMBLE

- Example of 2 trees
 - Prediction scores are summed to get the final score
 - Note that the two trees complement each other



SOME COMMENTS ABOUT ENSEMBLE TECHNIQUES

- Non-trivial increases in computational complexity for (usually) modest improvements in predictive performance
- \$1 million Netflix Grand Prize
 - For 3 years straight, thousands of teams could only improve upon the previous years results
 - Model tuning, factorization, simple ensembles etc
 - Could not beat the challenge of 10% improvement over the native Netflix recommendation algorithm
 - In 2009, a combined team (BellKors Pragmatic Chaos) beat the challenge using a complex blend of ensemble techniques
 - However, Netflix never implemented the winning model due to its complexity (development, execution and maintenance efforts that are not feasible for real-world application)

"[...] additional accuracy gains that we measured did not seem to justify the engineering effort needed to bring them into a production environment."

EXERCISE

- Upload “machine learning I.ipynb” into your Jupyter dashboard
 - Follow the comments and codes and feel free to try out new codes on your own
 - The visualization codes are meant to create graphs to help you picture the shape of the solution, don’t crack your head trying to understand the codes!
 - The point is to help you build your intuition to what the ML model does

UNSUPERVISED LEARNING

$$L(\mathbf{w}) = \prod_{i=1}^n P(y^i|x^i; \mathbf{w})$$

$$I_H(t) = -\sum_{i=1}^c p(i|t) \log_2 p(i|t)$$

$$k(x^i, x') = \exp\left(-\gamma \|x^i - x'\|^2\right)$$

$$p(Y \geq k) = \sum_{j=k}^n \binom{n}{j} \beta^j (1-\beta)^{n-j}$$

$$\varphi_2 = (u_2, v_2)$$

$$\varphi_3 = (u_3, v_3)$$

$$\varphi_4 = (u_4, v_4)$$

$$\varphi_5 = (u_5, v_5)$$

$$\varphi_6 = (u_6, v_6)$$

$$\varphi_7 = (u_7, v_7)$$

$$\varphi_8 = (u_8, v_8)$$

$$\varphi_9 = (u_9, v_9)$$

$$\varphi_{10} = (u_{10}, v_{10})$$

$$\varphi_{11} = (u_{11}, v_{11})$$

$$\varphi_{12} = (u_{12}, v_{12})$$

$$\varphi_{13} = (u_{13}, v_{13})$$

$$\varphi_{14} = (u_{14}, v_{14})$$

$$\varphi_{15} = (u_{15}, v_{15})$$

$$\varphi_{16} = (u_{16}, v_{16})$$

$$\varphi_{17} = (u_{17}, v_{17})$$

$$\varphi_{18} = (u_{18}, v_{18})$$

$$\varphi_{19} = (u_{19}, v_{19})$$

$$\varphi_{20} = (u_{20}, v_{20})$$

$$\varphi_{21} = (u_{21}, v_{21})$$

$$\varphi_{22} = (u_{22}, v_{22})$$

$$\varphi_{23} = (u_{23}, v_{23})$$

$$\varphi_{24} = (u_{24}, v_{24})$$

$$\varphi_{25} = (u_{25}, v_{25})$$

$$\varphi_{26} = (u_{26}, v_{26})$$

$$\varphi_{27} = (u_{27}, v_{27})$$

$$\varphi_{28} = (u_{28}, v_{28})$$

$$\varphi_{29} = (u_{29}, v_{29})$$

$$\varphi_{30} = (u_{30}, v_{30})$$

$$\varphi_{31} = (u_{31}, v_{31})$$

$$\varphi_{32} = (u_{32}, v_{32})$$

$$\varphi_{33} = (u_{33}, v_{33})$$

$$\varphi_{34} = (u_{34}, v_{34})$$

$$\varphi_{35} = (u_{35}, v_{35})$$

$$\varphi_{36} = (u_{36}, v_{36})$$

$$\varphi_{37} = (u_{37}, v_{37})$$

$$\varphi_{38} = (u_{38}, v_{38})$$

$$\varphi_{39} = (u_{39}, v_{39})$$

$$\varphi_{40} = (u_{40}, v_{40})$$

$$\varphi_{41} = (u_{41}, v_{41})$$

$$\varphi_{42} = (u_{42}, v_{42})$$

$$\varphi_{43} = (u_{43}, v_{43})$$

$$\varphi_{44} = (u_{44}, v_{44})$$

$$\varphi_{45} = (u_{45}, v_{45})$$

$$\varphi_{46} = (u_{46}, v_{46})$$

$$\varphi_{47} = (u_{47}, v_{47})$$

$$\varphi_{48} = (u_{48}, v_{48})$$

$$\varphi_{49} = (u_{49}, v_{49})$$

$$\varphi_{50} = (u_{50}, v_{50})$$

$$\varphi_{51} = (u_{51}, v_{51})$$

$$\varphi_{52} = (u_{52}, v_{52})$$

$$\varphi_{53} = (u_{53}, v_{53})$$

$$\varphi_{54} = (u_{54}, v_{54})$$

$$\varphi_{55} = (u_{55}, v_{55})$$

$$\varphi_{56} = (u_{56}, v_{56})$$

$$\varphi_{57} = (u_{57}, v_{57})$$

$$\varphi_{58} = (u_{58}, v_{58})$$

$$\varphi_{59} = (u_{59}, v_{59})$$

$$\varphi_{60} = (u_{60}, v_{60})$$

$$\varphi_{61} = (u_{61}, v_{61})$$

$$\varphi_{62} = (u_{62}, v_{62})$$

$$\varphi_{63} = (u_{63}, v_{63})$$

$$\varphi_{64} = (u_{64}, v_{64})$$

$$\varphi_{65} = (u_{65}, v_{65})$$

$$\varphi_{66} = (u_{66}, v_{66})$$

$$\varphi_{67} = (u_{67}, v_{67})$$

$$\varphi_{68} = (u_{68}, v_{68})$$

$$\varphi_{69} = (u_{69}, v_{69})$$

$$\varphi_{70} = (u_{70}, v_{70})$$

$$\varphi_{71} = (u_{71}, v_{71})$$

$$\varphi_{72} = (u_{72}, v_{72})$$

$$\varphi_{73} = (u_{73}, v_{73})$$

$$\varphi_{74} = (u_{74}, v_{74})$$

$$\varphi_{75} = (u_{75}, v_{75})$$

$$\varphi_{76} = (u_{76}, v_{76})$$

$$\varphi_{77} = (u_{77}, v_{77})$$

$$\varphi_{78} = (u_{78}, v_{78})$$

$$\varphi_{79} = (u_{79}, v_{79})$$

$$\varphi_{80} = (u_{80}, v_{80})$$

$$\varphi_{81} = (u_{81}, v_{81})$$

$$\varphi_{82} = (u_{82}, v_{82})$$

$$\varphi_{83} = (u_{83}, v_{83})$$

$$\varphi_{84} = (u_{84}, v_{84})$$

$$\varphi_{85} = (u_{85}, v_{85})$$

$$\varphi_{86} = (u_{86}, v_{86})$$

$$\varphi_{87} = (u_{87}, v_{87})$$

$$\varphi_{88} = (u_{88}, v_{88})$$

UNSUPERVISED LEARNING

- Up till this point, we have been dealing with problems where we have an answer to (outcome)
 - All our regression equations are setup with an outcome variable on the LHS
- What happens when we suspect that there are relationships, subsets or clusters in our data, but have no answers upfront?

CLUSTERING

- Clustering is a category of unsupervised learning that helps us discover hidden structures and natural grouping in data based on a **similarity function**
 - K-means
 - Hierarchical
 - Density-based
 - Graph-based (basis of network theory)
- Other unsupervised learning algorithms
 - Association rules (Apriori)
 - Mixture models, Factor models, and Latent models

K-MEANS CLUSTERING

- The most widely used clustering algorithm in academia and industry
 - Easy to implement
 - Computationally efficient
- Concept
 - Find groups of similar objects based on a distance function that makes them more related to each other than to objects in other groups
- Issues
 - Requires us to specify number of clusters k
 - Inappropriate choice of k can result in poor clustering
 - We can determine the optimal number of k using elbow and silhouette plots
 - Clusters do not overlap and are not hierarchical
 - Clusters can be empty!

STEPS TO IMPLEMENT K-MEANS

- Randomly pick k centroids from the sample points
 - Remember that we have to manually specify the number of clusters k
- Assign every sample to its nearest centroid
- Move the centroids to the center of the samples
- Repeat steps 2 and 3 until
 - Cluster assignments no longer change
 - User defined threshold has been reached
 - Number of iterations
 - Error rate

DISTANCE MEASURE

- The key mathematical function in clustering is the distance function
 - This is a measure of similarity between objects
 - The nearer 2 objects are together, than to others, the more likely they belong together in a cluster
- The most commonly used distance function is the squared Euclidean distance
 - Note that in 2-d form, it's the same as how we measure straight-line distance between 2 objects
 - The generalized vector form is: $d(\mathbf{x}, \mathbf{y})^2 = \sum_{j=1}^m (x_j - y_j)^2 = \|\mathbf{x} - \mathbf{y}\|_2^2$
- The objective is to minimize the within-cluster sum of squared errors

EXERCISE

- Upload “Machine Learning II.ipynb” into your dashboard
 - Follow the comments and instructions
 - Use the visualization to build your intuition as to the shape of the solution

MODEL EVALUATION AND HYPERPARAMETER TUNING

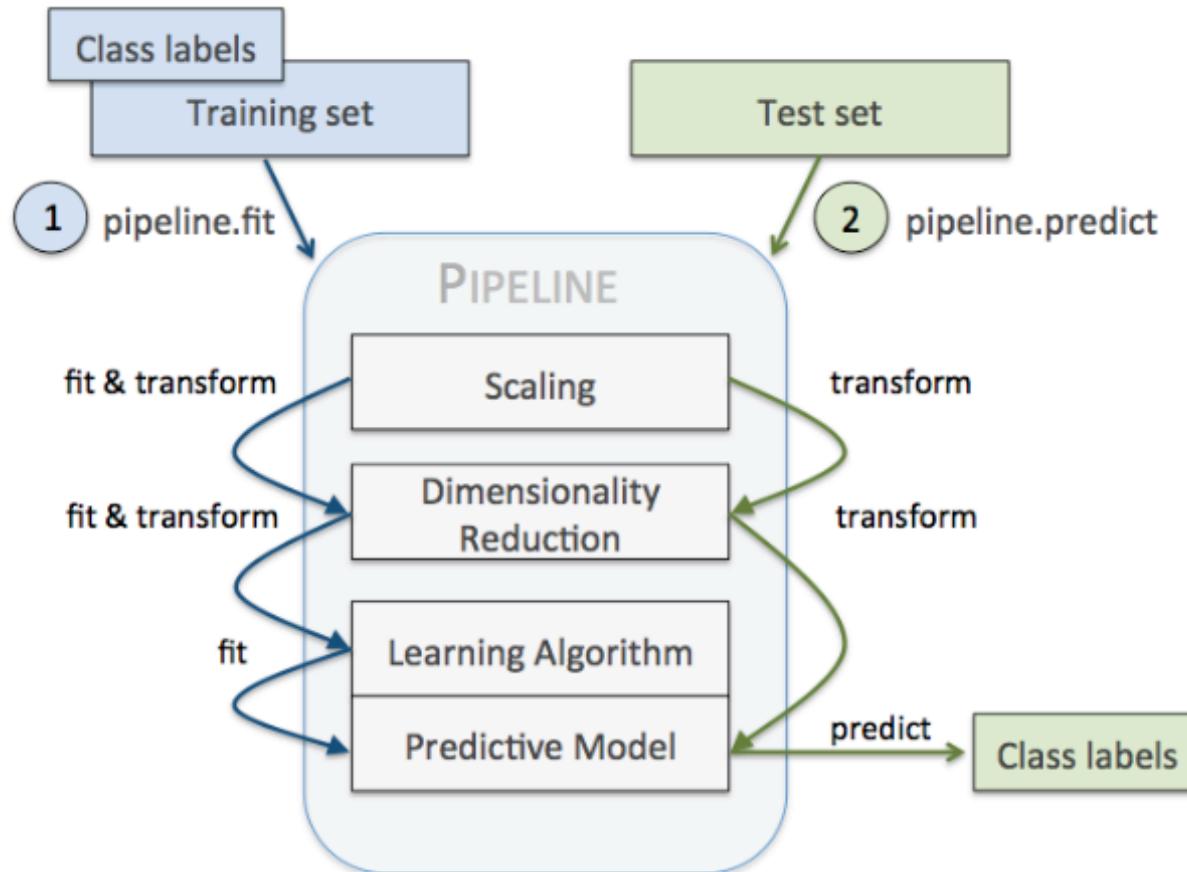
NEXT STEPS

- Previous sections have familiarized us with the basics of machine learning algorithms
- Building good models, however, require more than just knowing when to use which algorithms
- We need to have a framework for evaluating model performances in response to algorithm fine-tuning
 - Some of these are best practices contributed by the data science community

USING PIPELINES TO STREAMLINE WORKFLOWS

- You may have noticed that we have been using and reusing multiple data preprocessing, compression and hyperparameters
 - Some of these processes require that we fit training and testing datasets with the same parameters (e.g. scaling and compression transformation)
- Scikit-learn comes with a very useful tool that enhances this workflow
 - Known as the Pipeline

PIPELINE WORKFLOW

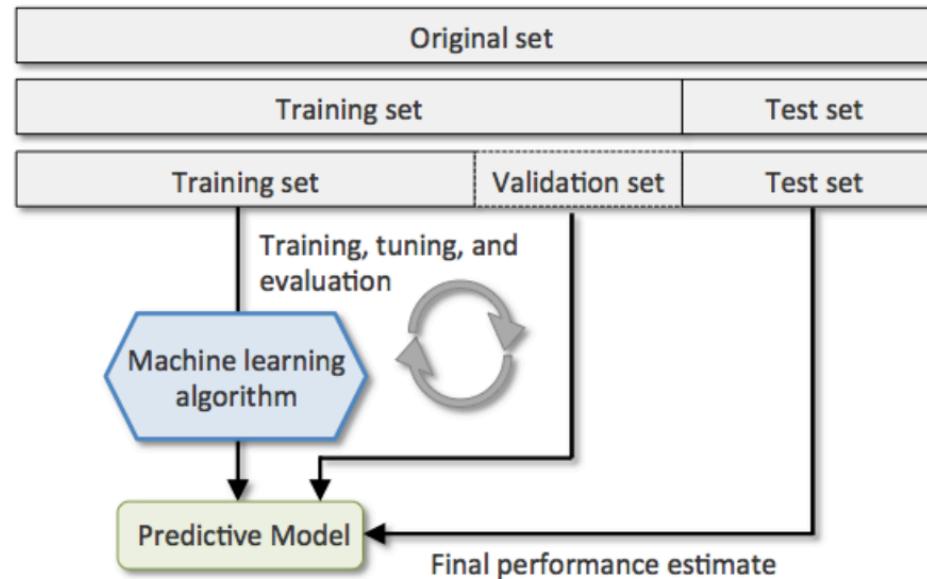


ASSESSING MODEL PERFORMANCE

- As covered in the previous sessions, model performance should be evaluated on 'never-seen-before' data or out-of-sample data
- Model performance can fall under the following 3 types:
 - Underfitting (model too simple, poor in-sample and out-sample performance)
 - Overfitting (model too complex, not generalizable and thus poor out-sample performance)
 - Good balance (this is what we want)
- We can use cross-validation techniques to figure out what the good balance is
 - Holdout method
 - K-fold

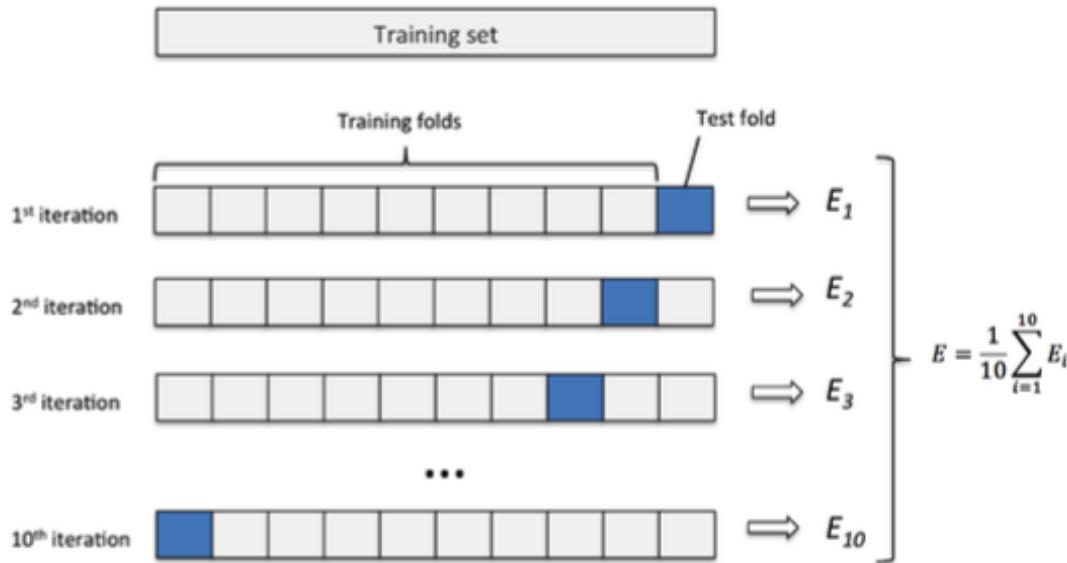
THE HOLDOUT METHOD

- Split the data into 3 sets:
 - Training set: used to fit the models
 - Validation set: models with the best performance on this set will be selected
 - Test set: final performance estimate
- However, performance estimates are sensitive to how we partition the data
 - Not the best CV method, though many still use this



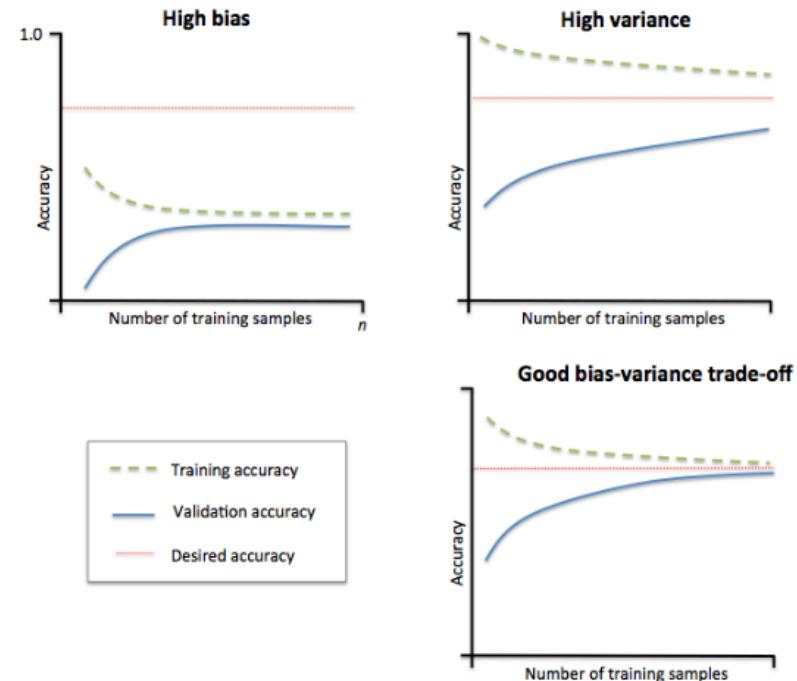
K-FOLD CROSS VALIDATION

- The preferred method for performance evaluation
 - 10 folds is the standard
 - But we will need to increase the number of folds for smaller and smaller datasets (up to leave-one-out for very small datasets)
 - Note that this will increase computational time
 - Also, the training samples will be pretty much similar to each other (the fitted models don't differ much from each other)
- To handle unequal class proportions (which can be a huge issue if the inequality is large), we can preserve the proportions in each fold
 - This is known as the stratified k-fold cross-validation



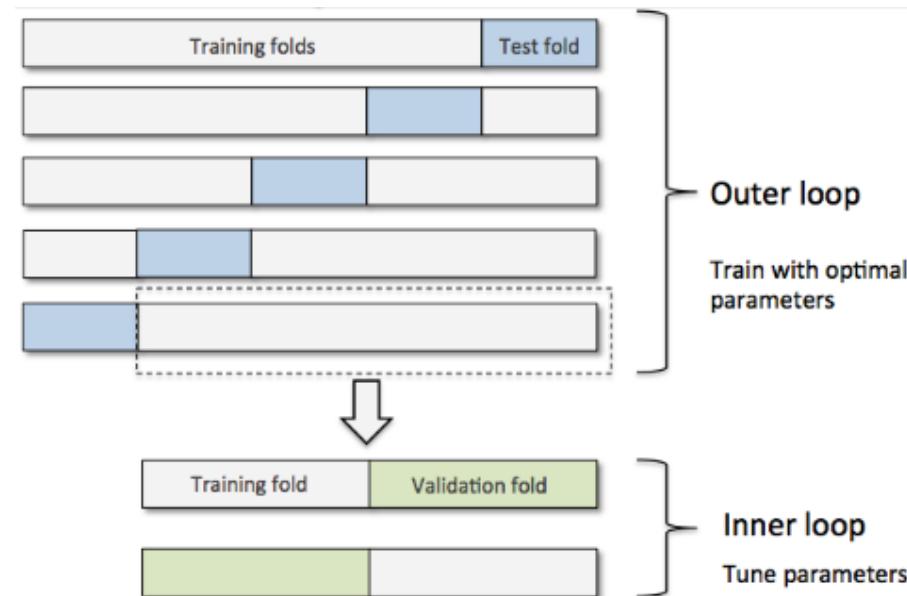
PERFORMANCE DIAGNOSTICS

- We can use learning and validation curves to diagnose our results
- Learning curves
 - Estimate underfitting vs overfitting
 - Graph on upper-left has low training and cv accuracy (underfitting)
 - Increase variables through collection or transformation
 - Decrease regularization (increasing c)
 - Graph on upper-right has high training but low cv accuracy (overfitting)
 - Increase observations
 - Decrease complexity through feature selection or extraction
 - Increase regularization (decreasing c)
- Validation curves
 - Improve models by varying parameters



OTHER USEFUL TECHNIQUES

- Grid search
 - Similar to validation curves, we specify a list of values for each hyperparameter
 - This method finds the optimal combination of hyperparameter values in a brute force exhaustive manner
- Nested cross-validation
 - Performing this test on different ML algorithms can aid us in choosing the best candidate model
 - Exhibits good performance in academic papers



CONFUSION MATRIX

- We can drill deeper into the accuracy measure to evaluate a model's precision and recall
- Precision, recall and the associated F1-score is represented by a confusion matrix
 - The % of true and false positive/negative
 - Precision is the % of true positives over true and false positives
 - How many of the predicted positives are indeed true?
 - Recall is the % of true positives over true positives + false negatives
 - How many of the actual positives are correctly predicted?
 - The F1 score is $2 \times (\text{Precision} \times \text{Recall}) / (\text{Precision} + \text{Recall})$

		Predicted class	
		P	N
		True Positives (TP)	False Negatives (FN)
Actual Class	P	True Positives (TP)	False Negatives (FN)
	N	False Positives (FP)	True Negatives (TN)

EXERCISE

- Upload “Tuning and Performance.ipynb” and “wdbc.csv” into your dashboard
 - Follow the comments and instructions
- As part of a final exercise, we will let you explore a financial dataset using the techniques that we have learnt
 - Upload “key_financials.xlsx” into your Jupyter dashboard
 - Upload “financial_exercise.ipynb” into your Jupyter dashboard
 - Follow the comments and codes and attempt the questions

QUESTIONS?

Email any queries to
jackhong@smu.edu.sg